

REAL-TIME MEASUREMENT OF MENTAL WORKLOAD: A FEASIBILITY STUDY

Arthur Kramer, Darryl Humphrey, Erik Sirevaag & Axel Mecklinger

Department of Psychology
University of Illinois

The primary goal of our study was to explore the utility of event-related brain potentials (ERP) as real-time measures of workload. To this end, subjects performed two different tasks both separately and together. One task required that subjects monitor a bank of constantly changing gauges and detect critical deviations. Difficulty was varied by changing the predictability of the gauges. The second task was mental arithmetic. Difficulty was varied by requiring subjects to perform operations on either two or three columns of numbers. Two conditions that could easily be distinguished on the basis of performance measures were selected for the real-time evaluation of ERPs. A bootstrapping approach was adopted in which one thousand samples of n trials ($n = 1, 3, 5 \dots 65$) were classified using several measures of P300 and Slow Wave amplitude. Classification accuracies of 85% were achieved with 25 trials. Results are discussed in terms of potential enhancements for real-time recording.

INTRODUCTION

The research presented here derives from an extensive series of investigations that have demonstrated the utility of Event-Related Brain Potentials (ERPs) in the assessment of residual capacity during the acquisition and performance of a variety of perceptual-motor and cognitive tasks (Donchin et al., 1986; Kramer, 1987). The focus of the present study was to assess the feasibility of employing ERPs as on-line measures of mental workload. If physiological data, and ERPs in particular, are to serve as real-time measures of operator mental load, the amount of data (e.g. secs, mins?) necessary to reliably discriminate among levels of workload must be determined. This question will be addressed in the present study by adopting a bootstrapping approach in which we examine the classification accuracy of ERP measures with from 1 to 65 secs of data. However, before we describe our experiment in detail we will briefly discuss the previous research that suggests that ERPs provide a sensitive and reliable measure of mental load in an off-line context.

Several recent studies have illustrated the usefulness of the ERP,

and more specifically the P300 component, as an index of processing resources (Horst et al., 1984; Isreal et al., 1980; Kramer et al., 1985, 1987; Natani and Gomer, 1981; Sirevaag et al., 1988). The general paradigm employed in these studies requires subjects to perform two tasks concurrently. One task is designated as primary and the other task as secondary. Subjects are instructed to maximize their performance on the primary task and devote any additional resources to the performance of the secondary task.

ERPs are elicited by events in either one or both of the tasks. Increases in the perceptual/cognitive difficulty of the primary task result in a decrease in the amplitude of the P300s elicited by the secondary task. Conversely, P300s elicited by discrete events embedded within the primary task increase in amplitude with increases in primary task difficulty. Furthermore, changes in response related demands of a task have no influence on the P300 (Isreal et al., 1980).

The reciprocal relationship between P300s elicited by primary and secondary task stimuli is consistent with the resource tradeoffs presumed to underlie dual-task performance decrements (Kahneman, 1973; Navon and Gopher, 1979; Sanders, 1979; Wickens, 1980). That is, resource models predict that as the difficulty of one task is increased, additional resources are re-allocated to that task in order to maintain performance, thereby depleting the supply of resources that could have been used in the processing of other tasks. Thus, the P300 appears to provide a measure of resource tradeoffs that can only be inferred from more traditional performance measures. Furthermore, P300s elicited by secondary task events are selectively sensitive to the perceptual/cognitive demands imposed upon the operator. This selective sensitivity may be especially useful in decomposing the changing processing requirements of complex tasks (Kramer, 1987).

One might ask why ERPs should be used to monitor changes in resource demands given that several technically simpler approaches to the assessment of skill acquisition and mental workload have already been implemented. Although

numerous performance-based measures of mental workload exist, they suffer from several drawbacks. First, some of the measurement techniques require subjects to perform a secondary task which frequently interferes with the performance of the task of interest (Knowles, 1963; Rolfe, 1971; Wickens, 1979). This is clearly unacceptable in an operational environment in which the safety of the operator must be assured. Even in the laboratory setting it is difficult to determine which of the two tasks generated an observed performance decrement since the performance on the two tasks is easily confounded. Second, performance-based measures of mental workload provide an output measure of the operator's information processing activities (e.g. RT, accuracy). Thus, at best, performance measures provide only an indirect index of cognitive function. Third, performance measures do not always correlate highly with the actual workload of the tasks (Brown, 1978; Dornic, 1980; Ogden et al., 1979). Fourth, although subjective measures are relatively easy to collect and possess high face validity they do not reflect the moment to moment variations in workload that can be indexed by physiological measures.

The present study is part of a continuing effort to explore the utility of psychophysiological measures of mental workload. A primary aim of the project is to determine the feasibility of on-line uses of integrated psychophysiological and performance data. However, given the magnitude of the project this report will be confined to a description of a preliminary examination of signal/noise ratio parameters of ERPs. More specifically, we will derive the functions that relate amount of ERP data to discrimination accuracy between workload conditions.

METHODS

Subjects

Four dextral subjects (2 female) were paid \$4.00/hour plus a dollar/day bonus for their participation in five sessions. All subjects had normal or corrected-to-normal vision.

Tasks

Two different tasks were performed both separately and together. We will describe each of the tasks in detail.

Monitoring Task. One task consisted of monitoring six gauges. The behavior of a gauge was determined by the interaction of four properties: update speed, noise level, noise frequency and transients. The cursors moved around the gauges at different speeds, a slower gauge taking longer to

reach the critical region. Noise level was the amount of random jitter in the cursor. Noise frequency determined how often random fluctuations were added to a gauge. The addition of transients also served to perturb a gauge.

The interaction of these properties produced cursor driving functions of varying predictability. Manipulating the driving functions allowed control over gauge monitoring difficulty. The driving functions employed in the high predictability (HP) conditions were such that within a row of three gauges the driving functions were identical in terms of speed, noise level, and noise frequency; no transient occurred for any gauge. The two rows differed in the speed of cursor movement, speed being constant within a row. For the low predictability (LP) conditions the average value for all properties was equivalent to the HP conditions, however, the individual values were varied with no established correlation between any set of gauges. The LP conditions contained three gauges with a transient. The frequency of the transient was different for each of the three gauges.

The gauges were presented on a CRT in front of the subject. Each gauge was divided into 12 regions (labelled 1 to 12). In addition, each third of the gauge was distinctly colored (green, yellow and red). The critical level was designated by the position marked by the numeral 9, which was the first region in the red zone.

The purpose of this task was to reset each gauge as quickly as possible once its cursor had entered the critical region. To reset a gauge the subjects pressed one of six keys after which the cursor returned to the starting position marked by the numeral 1. The cursors were not continuously visible. To sample a given gauge the subject pressed one of a set of six keys with their left hand. The cursor remained visible for 1000 msec. Simultaneous sampling was not possible.

Mental Arithmetic Task. The center of each gauge served as a display area for the operands and operators of the mental arithmetic trials. All of the operands and operators were presented simultaneously and remained in view until an answer was entered or for a maximum of 30 seconds. An answer window appeared to the right of the gauges. Answers were entered via the numeric keypad of the response keyboard and appeared in the window as they were typed. Completion was signaled by pressing the 'enter' key of the numeric keypad. The inter-trial interval varied from four to fifteen seconds. Difficulty was manipulated by varying the number of

column operations necessary to complete the problem. The easy version of the task required operations on two columns while the difficult version of the task required operations on three columns of numbers. Henceforth, these versions of the tasks will be referred to as A2 and A3, respectively. Operations included addition and multiplication.

Subjects participated in five sessions. The first two sessions constituted training. Single task conditions, starting with the easy conditions progressing to the difficult conditions were performed first, followed by the dual task conditions. In the final three sessions the subject performed the eight conditions in a random order determined by a Latin square design. Only the data from the last three sessions will be presented in this report. In all sessions two blocks of each condition were run consecutively, each block taking five minutes. A five minute break was imposed at the halfway point in addition to any breaks the subject requested.

Performing the gauge monitoring and mental arithmetic tasks in all possible combinations yields eight conditions: 2 task types X 2 levels of difficulty X 2 task pairings (single or dual task condition).

ERP Recording

Electroencephalographic (EEG) activity was recorded from three midline sites (Fz, Cz, Pz according to the International 10-20 system: Jasper 1958) referenced to averaged mastoids. All electrodes were Sensormedics Ag/AgCl electrodes. The scalp electrodes were affixed with Grass EC2 electrode cream. The forehead ground, mastoid and electrooculogram (EOG) electrodes were affixed with the Grass cream and electrode collars. Vertical and horizontal EOG was in order to control for eye movement artifacts. Electrode impedance was maintained below 10 kohms.

The EEG and EOG were amplified by Grass 12A5 amplifiers with a 8 sec time constant and a low-pass filter of 100 Hz. The recording epoch was 1300 msec beginning 100 msec prior to an event. The data channels were digitized every 5 msec and were filtered off-line (-3 db at 6.89 Hz., 0 db at 22.22 Hz) prior to further analysis. The psychophysiological data collection was governed by a DEC PDP 11/73 computer system. Artifact rejection was based upon the vertical eye movement standard deviation. ERPs were recorded during the three experimental sessions.

Subjects were seated in a dimly lit, sound attenuated booth. Stimuli were presented on a color monitor located 80

cm in front of the subject. Stimulus presentation and behavioral data collection were performed by an IBM AT.

Data Analysis Procedures

ERP eliciting events included critical gauge samples, non-critical gauge samples and, presentation of math trials. ERP measurements included P300 latency, P300 base-to-peak amplitude, P300 base-to-peak area and, slow wave area. Behavioral variables included accuracy and response speed in both the monitoring and arithmetic tasks.

In an effort to determine the amount of physiological data needed to discriminate among different experimental conditions we applied a bootstrapping approach to single trial ERP data. Given the amount of data collected in our study we decided to begin by examining the physiological differences between two conditions that could be discriminated on the basis of performance measures: the LP single task gauge condition and the gauge samples from the LP/A3 dual task conditions. One thousand samples of size n ($n = 1, 3, 5, \dots, 65$) were randomly selected from single trial data in each of these conditions. By comparing the single trial samples with the grand average waveforms for that condition the single trial may be classified as a hit (belonging to the criterion condition), a miss (not belonging to the criterion condition) or unclassifiable. Tabulating the classification results in a 2 X 2 contingency table enabled us to assess the efficiency of a number of ERP measures.

RESULTS & DISCUSSION

The results will be organized in the following manner. First, we will describe the effects of single and dual task manipulations on subjects' performance and ERPs. These analyses will enable us to establish the relative differences in performance and workload among the single and dual task conditions. Second, we will select two experimental conditions that can be distinguished on the basis of average performance and ERP measures. A bootstrapping approach will then be applied to the single trial ERP data in these conditions. The classification accuracy value derived from each sample of one thousand measures will then be plotted as a function of the number of trials in each of the thousand samples. This procedure enables us to determine how changes in the signal/noise ratio of the ERP as a function of averaging (e.g. averaging from 1 to 65 trials for each of the thousand samples) translates into gains in the accuracy of discrimination between workload conditions.

The bootstrapping approach will be applied to several different ERP measures including: base to peak measures of P300 amplitude (P3bp), measures of P300 area (P3area), cross-correlation measures of P300 amplitude (P3cross), and area measures of a late slow wave component (SWarea). P3bp was defined as the largest positivity in the waveform between 300 and 800 msec post-stimulus relative to a pre-stimulus baseline. The "stimulus" was the presentation of the cursor with the gauges. P3area was defined as the area in a 100 msec window centered around the peak. P3cross measures were calculated by moving a 300 msec wide cosine wave across the period from 300 to 800 msec post-stimulus. The slope of the regression function at the point at which the correlation between the cosine "template" and the ERP waveform was maximized was defined as P3cross. SWarea was defined as the area between 750 and 1000 msec post-stimulus.

Effects of Experimental Manipulations

Figure 1 presents a measure of the accuracy with which subjects reset the gauges in each of the monitoring conditions. A "miss" was scored when subjects failed to reset a gauge within 10 sec following the point at which it reached a critical value. As can be seen from the figure, accuracy decreased from single to dual task conditions and again with an increase in the difficulty of the dual task. Accuracy also appeared to differ as a function of the predictability of the gauges (HP vs. LP). These differences were confirmed by a repeated measures 2-way ANOVA, with gauge (2 gauge conditions, HP and LP) and task (3 arithmetic conditions, none, A2 and A3) as factors. Significant main effects were obtained for both the gauge ($F(1,3)=13.2$, $p<.01$) and task ($F(2,6)=21.2$, $p<.01$) factors. A marginally significant interaction between gauge and task factors was also obtained ($F(2,6)=2.9$, $p<.08$) suggesting a decrease in accuracy at the most difficult level of each of the factors.

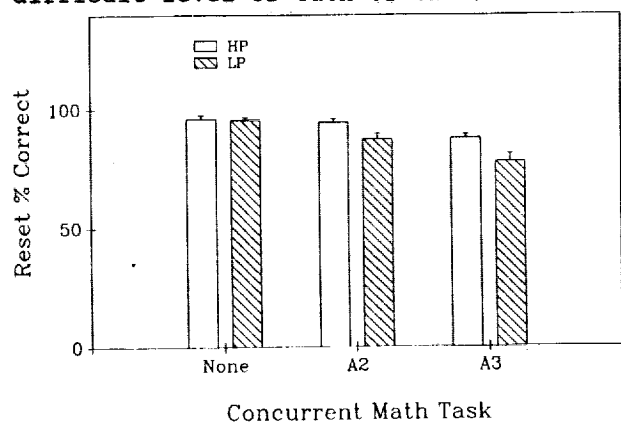


Figure 1. Accuracy in the monitoring task.

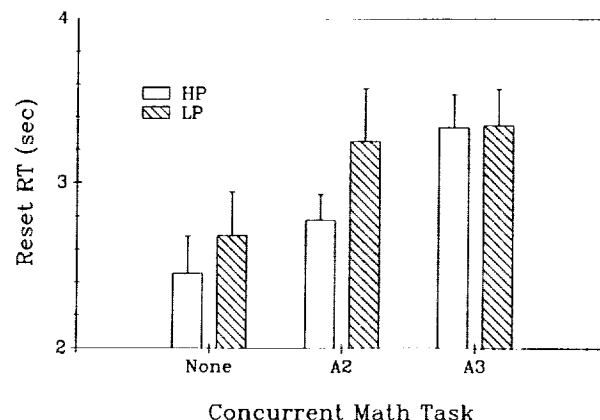


Figure 2. Reset RT in the monitoring task.

Figure 2 presents gauge reset RTs for each of the monitoring conditions. A repeated measures ANOVA performed on this data set revealed a significant main effect for the task factor ($F(2,6)=5.4$, $p<.01$). RT increased from the single to the dual task conditions and again from the A2 to the A3 versions of the arithmetic task. The main effect for the gauge factor did not attain statistical significance.

Accuracy and RT measures are presented for the arithmetic task in figures 3 and 4, respectively. Accuracy in the arithmetic task was higher when operations were performed on two columns than when a three column problem was performed ($F(1,3)=22.8$, $p<.01$). RT was also faster in the A2 than in the A3 version of the arithmetic task ($F(1,3)=26.4$, $p<.01$). Finally, RT in the arithmetic task increased with the transition from the single to dual task conditions and again when the difficulty of the monitoring task was increased.

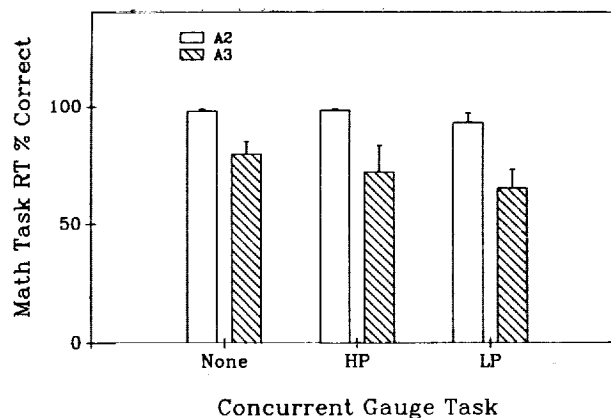


Figure 3. Accuracy in the arithmetic task.

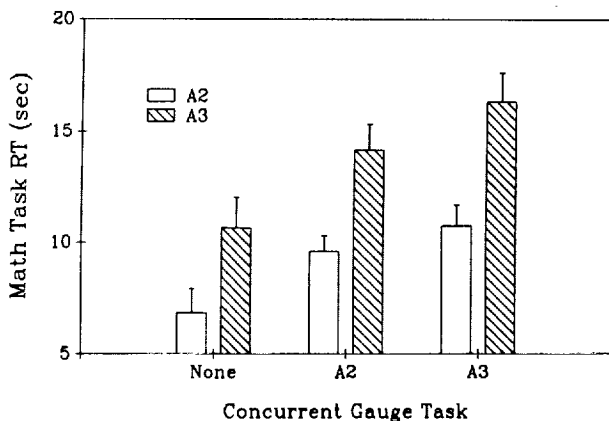


Figure 4. RT in the arithmetic task.

Real-Time Analysis of Mental Workload

Given the substantial amount of analysis time required to perform the "bootstrapping" operation we decided to select two experimental conditions to analyze further. In order to perform the bootstrapping operation it was necessary for the experimental conditions to meet three criteria. First, there should be a substantial number of trials available in the selected conditions. This was necessary since repeated samples of 1000 trials would be selected during the bootstrapping operation. Second, the conditions should be discriminable on the basis of performance measures. Thus, we wanted to begin our analysis of the real-time potential of ERPs by selecting two clearly discriminable conditions. Later analyses will examine conditions that are less discriminable. Third, the conditions should be discriminable on the basis of average ERP measures. Based on these criteria we selected two conditions from the monitoring task: the single task LP condition and the dual task LP/A3 condition.

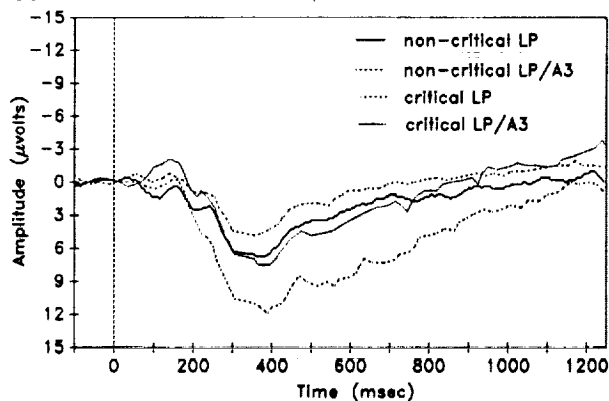


Figure 5. Grand average ERPs recorded at Pz for four of the monitoring conditions.

Figure 5 presents the grand average ERPs across the four subjects for the LP and LP/A3 conditions. It is important to note that we have further subdivided the conditions into waveforms that were elicited during times at which the gauges were in the acceptable range and other times in which the gauges had gone critical. Since the gauge critical samples were most closely associated with the performance measures we decided to employ ERPs to discriminate between the LP and LP/A3 conditions during the gauge critical periods. Approximately 200 trials were available in each of these conditions for each of the subjects. The bootstrapping operation was performed separately on the data from two of the original four subjects.

As described above, the bootstrapping operation involved the repeated selection of single trial ERPs from each of the conditions. Each "sample" was comprised of 1000 ERP measures, 500 selected from the LP condition and 500 selected from the LP/A3 condition. Each of the ERP measures was composed of an average of from 1 to 65 single trial ERP waveforms. Classification accuracy was determined by computing the relative "distance" of each ERP measure from the subject's grand average ERP measures in the LP and LP/A3 conditions. For example, if a subject possessed a grand average P300 amplitude of 50 microvolts in the LP/A3 condition and 10 microvolts in the LP condition then a single trial measure of 46 microvolts would be classified as LP/A3. This classification procedure was performed for each of the 1000 ERP measures in a sample and for each of the different pattern recognition techniques (i.e. P3bp, P3area, P3cross, SWcross).

Figures 6 and 7 present the classification functions for subjects 2 and 3, respectively. In the figures we plot the accuracy of classification (y-axis) against the number of single trial ERPs that were averaged to produce each of the ERP measures in a sample (each sample included 1000 ERP measures). Several aspects of the figures are noteworthy. First, for each of the pattern recognition techniques plotted, classification accuracy increased with increases in the number of trials per measure. This continued improvement in classification accuracy represents the increasing signal/noise ratio as additional single trials are averaged to produce each measure. Second, it is clear from the figures that the pattern recognition techniques improved at different rates and achieved different asymptotic levels of accuracy. For both of the subjects P3bp and P3area improved more quickly and achieved higher levels of performance than SWarea and P3cross. In fact, P3cross is not plotted for subject 2 because it never exceeded 50% classification accuracy.

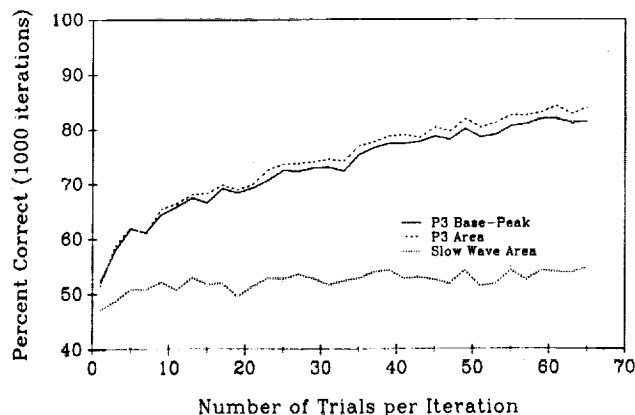


Figure 6. Classification accuracy as a function of the number of trials per measure for subject 2.

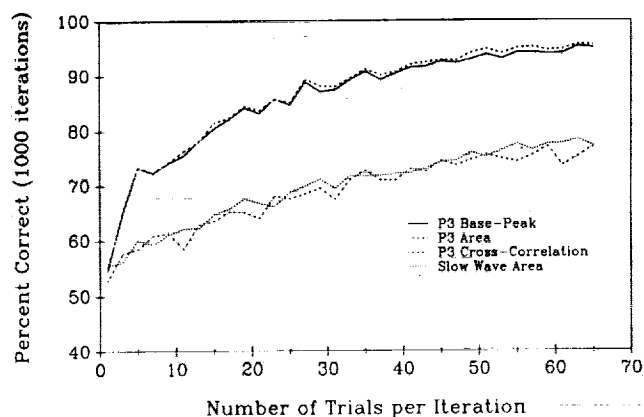


Figure 7. Classification accuracy as a function of the number of trials per measure for subject 3.

Third, for both P3bp and P3area there was a dramatic improvement in classification accuracy with the addition of the first five single trials followed by a more gradual improvement as additional trials were averaged. Finally, it is interesting to note that classification accuracy improved and reached different asymptotic levels for the two subjects.

SUMMARY AND CONCLUSIONS

The results of our investigation provide support for the utility of ERPs as real-time measures of mental workload. However, it is important to note that this support is both preliminary and tentative due to the small number of subjects, conditions, and pattern recognition techniques used in our study. The results are encouraging, however, and suggest a number of avenues for further exploration.

First, the differential efficiency of the pattern recognition techniques suggests that other techniques may offer improvements over the four that we have

examined. In our study we used techniques that capitalized on the differences between only one component of the ERP (i.e. either P300 or Slow Wave amplitude). However, a number of other ERP components also appear to be sensitive to variations in mental workload (Horst et al., 1984; Kramer, 1987). Given that these components reflect changes in workload not indexed by P300 and Slow Wave amplitude, the use of multivariate techniques such as discriminant functions should improve the ability to discriminate among different levels of workload. It might also be possible to enhance discriminability by examining changes in the frequency spectra of EEG.

Second, previous examinations of the accuracy of single trial classifications of ERPs have suggested that the efficiency of different pattern recognition techniques is dependent on the characteristics of subject's waveforms (Farwell and Donchin, 1988). For example, base to peak measures tend to be most successful when the component of interest is sharply defined while area measures are superior for wider components. Differences in the efficiency of P3cross and SWarea measures for our two subjects also appear to be due to differences in their waveforms. Thus, these analyses suggest that it might be useful to compile a set of heuristics that map waveform characteristics to pattern recognition techniques.

Third, it seems reasonable to suppose that the ability to discriminate among workload levels depends on the homogeneity within workload levels. In the present study we selected gauge samples in the LP/A3 condition irrespective of whether subjects were performing the arithmetic task (arithmetic tasks were presented with isi's of from 4 to 15 secs). Thus, our LP/A3 condition was actually a mixture of single and dual task trials. A comparison of the "dual task" trials in the LP/A3 condition with the LP condition should increase classification accuracy.

Fourth, while it is important to determine classification accuracy in the "best-case" situation it is also imperative that classification functions are derived for smaller differences in workload. We are currently examining the range of sensitivity of ERP measures to graded differences in workload. Finally, it is clear that classification accuracy can be improved by integrating psychophysiological and performance measures into predicative and descriptive equations. Therefore, it is necessary to determine how the relative sensitivity of different physiological and performance measures vary with changes in task structure and subject

state.

ACKNOWLEDGMENT

The research was supported by a grant from the Office of Naval Research (N00014-89-J-1493) with Dr. Harold Hawkins as technical monitor and by a grant from NASA Ames Research Center (NASA NAG-2-308) monitored by Dr. Sandra Hart.

REFERENCES

Brown, I. (1978). Dual-task methods of assessing workload. Ergonomics, 21, 221-224.

Donchin, E., Kramer, A., & Wickens, C. (1986). Applications of brain event-related potentials to problems in engineering psychology. In M. Coles, S. Porges & E. Donchin (Eds.), Psychophysiology: Systems, processes and applications. New York: Guilford.

Dornic, S. (1980). Language dominance, spare capacity and perceived effort in bilinguals. Ergonomics, 23, 369-377.

Farwell, L. & Donchin, E. (1988). Talking off the top of your head: toward a mental prosthesis using event-related brain potentials. Electroencephalography and Clinical Neurophysiology, 70, 510-523.

Horst, R., Munson, R. & Ruchkin, D. (1984). Event related brain potential indices of workload in a single task paradigm. Proceedings of the Human Factors Society, 28th Annual Meeting. San Antonio, Texas.

Isreal, J., Chesney, J., Wickens, C. & Donchin, E. (1980). P300 and tracking difficulty: Evidence for multiple resources in dual task performance. Psychophysiology, 17, 259-273.

Jasper, H. (1958). The ten-twenty electrode system of the international federation. Electroencephalography and Clinical Neurophysiology, 10, 371-375.

Kahneman, D. (1973). Attention and effort. Englewood Cliffs, N.J.: Prentice Hall.

Knowles, W. (1963). Operator loading tasks. Human Factors, 5, 155-161.

Kramer, A. (1987). Event-related brain potentials. In A. Gale & B. Christie (Eds.), Psychophysiology and the electronic workplace. Sussex, England: Wiley.

Kramer, A., Sirevaag, E. & Braune, R. (1987). A psychophysiological assessment of operator workload during

simulated flight missions. Human Factors, 29, 145-160.

Kramer, A., Wickens, C. & Donchin, E. (1985). The processing of stimulus attributes: Evidence for dual-task integrality. Journal of Experimental Psychology: Human Perception and Performance, 11, 393-408.

Natani, K. & Gomer, F. (1981). Electrocortical activity and operator workload: A comparison of changes in the electroencephalogram and in event related potentials. McDonnell Douglas Technical Report, MDC E2427.

Navon, D. & Gopher, D. (1979). On the economy of the human information processing system. Psychological Review, 86, 214-255.

Ogden, G., Levine, J. & Eisner, E. (1979). Measurement of workload by secondary tasks. Human Factors, 21, 529-548.

Rolfe, J. (1971). The secondary task as a measure of mental load. In W. Singleton, J. Fox & D. Witfield (Eds.), Measurements of man at work. London: Taylor and Francis.

Sanders, A. (1979). Some remarks on mental load. In N. Moray (Ed.), Mental workload: Its theory and measurement. New York: Plenum Press.

Sirevaag, E., Kramer, A., Coles, M. & Donchin, E. (1988). Resource reciprocity: An event-related brain potentials analysis. Acta Psychologica, 70, 77-97.

Wickens, C. (1979). Measures of workload, stress and secondary tasks. In N. Moray (Ed.), Mental workload: Its theory and measurement. New York: Plenum Press.

Wickens, C. (1980). The structure of attentional resources. In R. Nickerson & R. Pew (Eds.), Attention and Performance VIII. Hillsdale, N.J.: Erlbaum.

